

# CSCI567 Machine Learning (Fall 2018)

Prof. Haipeng Luo

U of Southern California

Sep 26, 2018

# Administration

HW 2 released, due on 10/07

HW 1 solutions are available, W1 is graded, P1 grading should be done by the end of this week

# Administration

HW 2 released, due on 10/07

HW 1 solutions are available, W1 is graded, P1 grading should be done by the end of this week

## Midterm

- 10/03, 5PM-7PM, *THH 101 and THH 201* (pick any one), sit every other seat, closed-book/notes, no electronics allowed
- Multiple choice + regular problems like homework
- Coverage: mostly Lec 1-5, some simple questions from Lec 6 (today)

# Outline

- 1 Review of last lecture
- 2 Support vector machines (primal formulation)
- 3 A detour of Lagrangian duality
- 4 Support vector machines (dual formulation)

# Outline

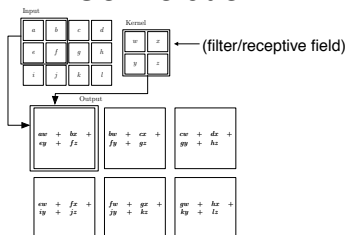
- 1 Review of last lecture
- 2 Support vector machines (primal formulation)
- 3 A detour of Lagrangian duality
- 4 Support vector machines (dual formulation)

# Convolutional Neural Nets

## Typical architecture for CNNs:

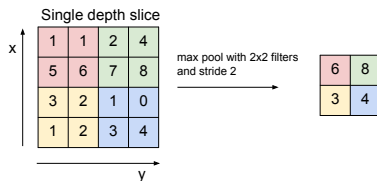
Input  $\rightarrow$  [[Conv  $\rightarrow$  ReLU]\*N  $\rightarrow$  Pool?]\*M  $\rightarrow$  [FC  $\rightarrow$  ReLU]\*Q  $\rightarrow$  FC

## 2D Convolution



(Goodfellow 2016)

## MAX POOLING



Fei-Fei Li &amp; Justin Johnson &amp; Serena Yeung

Lecture 5 - 73

April 18, 2017

# Kernel functions

**Definition:** a function  $k : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$  is called a *(positive semidefinite) kernel function* if there exists a function  $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$  so that for any  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^D$ ,

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

# Kernel functions

**Definition:** a function  $k : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$  is called a *(positive semidefinite) kernel function* if there exists a function  $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$  so that for any  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^D$ ,

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

Examples we have seen

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^2$$

$$k(\mathbf{x}, \mathbf{x}') = \sum_{d=1}^D \frac{\sin(2\pi(x_d - x'_d))}{x_d - x'_d}$$

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^d \quad \text{(polynomial kernel)}$$

$$k(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\sigma^2}} \quad \text{(Gaussian/RBF kernel)}$$



# Kernelizing ML algorithms

Feasible as long as **only inner products are required**:

- regularized linear regression (dual formulation)

$$\phi(\mathbf{x})^T \mathbf{w}^* = \phi(\mathbf{x})^T \Phi^T (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y} \quad (\mathbf{K} = \Phi \Phi^T \text{ is } \textit{kernel matrix})$$

- nearest neighbor classifier with L2 distance

$$\|\phi(\mathbf{x}) - \phi(\mathbf{x}')\|_2^2 = k(\mathbf{x}, \mathbf{x}) + k(\mathbf{x}', \mathbf{x}') - 2k(\mathbf{x}, \mathbf{x}')$$

- perceptron, logistic regression, SVM, ...

# Outline

- 1 Review of last lecture
- 2 Support vector machines (primal formulation)**
- 3 A detour of Lagrangian duality
- 4 Support vector machines (dual formulation)

# Support vector machines (SVM)

- One of the most commonly used classification algorithms
- Works well with the kernel trick
- Strong theoretical guarantees

# Support vector machines (SVM)

- One of the most commonly used classification algorithms
- Works well with the kernel trick
- Strong theoretical guarantees

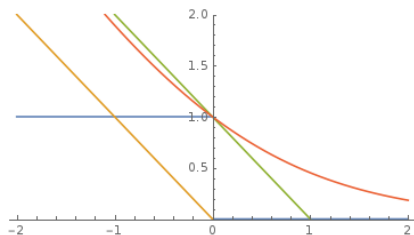
We focus on **binary classification** here.

# Primal formulation

In one sentence: linear model with L2 regularized hinge loss.

# Primal formulation

In one sentence: linear model with L2 regularized hinge loss. Recall



- **perceptron loss**  $\ell_{\text{perceptron}}(z) = \max\{0, -z\} \rightarrow$  Perceptron
- **logistic loss**  $\ell_{\text{logistic}}(z) = \log(1 + \exp(-z)) \rightarrow$  logistic regression
- **hinge loss**  $\ell_{\text{hinge}}(z) = \max\{0, 1 - z\} \rightarrow$  **SVM**

# Primal formulation

For a linear model  $(\mathbf{w}, b)$ , this means

$$\min_{\mathbf{w}, b} \sum_n \max \{0, 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)\} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

- recall  $y_n \in \{-1, +1\}$
- a nonlinear mapping  $\phi$  is applied
- the bias/intercept term  $b$  is used explicitly (think about why after this lecture)

# Primal formulation

For a linear model  $(\mathbf{w}, b)$ , this means

$$\min_{\mathbf{w}, b} \sum_n \max \{0, 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)\} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

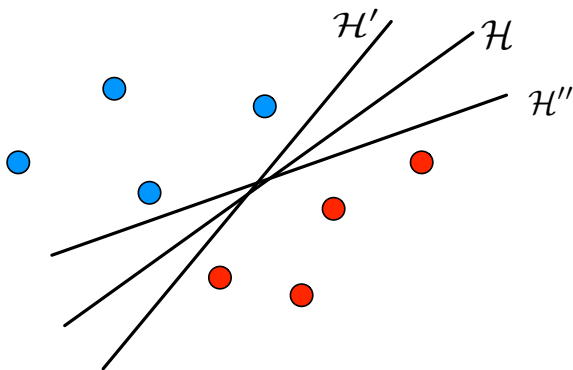
- recall  $y_n \in \{-1, +1\}$
- a nonlinear mapping  $\phi$  is applied
- the bias/intercept term  $b$  is used explicitly (think about why after this lecture)

*So why L2 regularized hinge loss?*



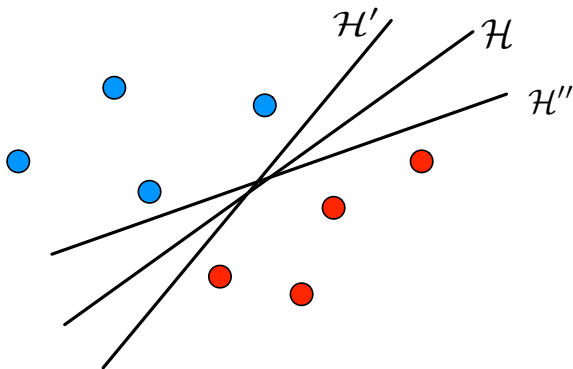
## Geometric motivation: separable case

When data is **linearly separable**, there are *infinitely many hyperplanes with zero training error*:



## Geometric motivation: separable case

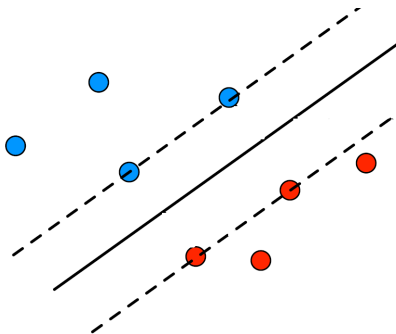
When data is **linearly separable**, there are *infinitely many hyperplanes with zero training error*:



So which one should we choose?

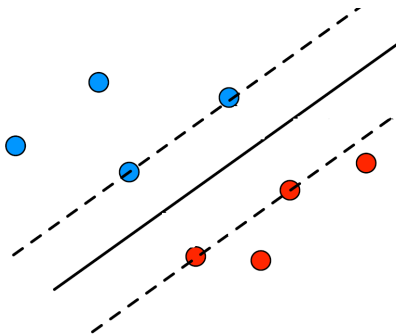
# Intuition

The further away from data points the better.



# Intuition

The further away from data points the better.



*How to formalize this intuition?*

## Distance to hyperplane

What is the **distance** from a point  $\mathbf{x}$  to a hyperplane  $\{\mathbf{x} : \mathbf{w}^T \mathbf{x} + b = 0\}$ ?

## Distance to hyperplane

What is the **distance** from a point  $\mathbf{x}$  to a hyperplane  $\{\mathbf{x} : \mathbf{w}^T \mathbf{x} + b = 0\}$ ?

Assume the **projection** is  $\mathbf{x} - \ell \frac{\mathbf{w}}{\|\mathbf{w}\|_2}$ , then

$$0 = \mathbf{w}^T \left( \mathbf{x} - \ell \frac{\mathbf{w}}{\|\mathbf{w}\|_2} \right) + b = \mathbf{w}^T \mathbf{x} - \ell \|\mathbf{w}\| + b$$

and thus  $\ell = \frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|_2}$ .

## Distance to hyperplane

What is the **distance** from a point  $\mathbf{x}$  to a hyperplane  $\{\mathbf{x} : \mathbf{w}^T \mathbf{x} + b = 0\}$ ?

Assume the **projection** is  $\mathbf{x} - \ell \frac{\mathbf{w}}{\|\mathbf{w}\|_2}$ , then

$$0 = \mathbf{w}^T \left( \mathbf{x} - \ell \frac{\mathbf{w}}{\|\mathbf{w}\|_2} \right) + b = \mathbf{w}^T \mathbf{x} - \ell \|\mathbf{w}\| + b$$

and thus  $\ell = \frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|_2}$ .

Therefore the distance is

$$\frac{|\mathbf{w}^T \mathbf{x} + b|}{\|\mathbf{w}\|_2}$$

## Distance to hyperplane

What is the **distance** from a point  $\mathbf{x}$  to a hyperplane  $\{\mathbf{x} : \mathbf{w}^T \mathbf{x} + b = 0\}$ ?

Assume the **projection** is  $\mathbf{x} - \ell \frac{\mathbf{w}}{\|\mathbf{w}\|_2}$ , then

$$0 = \mathbf{w}^T \left( \mathbf{x} - \ell \frac{\mathbf{w}}{\|\mathbf{w}\|_2} \right) + b = \mathbf{w}^T \mathbf{x} - \ell \|\mathbf{w}\| + b$$

and thus  $\ell = \frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|_2}$ .

Therefore the distance is

$$\frac{|\mathbf{w}^T \mathbf{x} + b|}{\|\mathbf{w}\|_2}$$

For a hyperplane that correctly classifies  $(\mathbf{x}, y)$ , the distance becomes

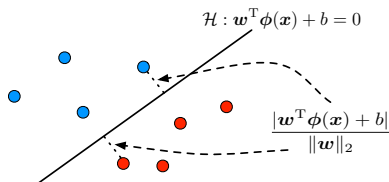
$$\frac{y(\mathbf{w}^T \mathbf{x} + b)}{\|\mathbf{w}\|_2}$$



# Maximizing margin

**Margin:** the *smallest* distance from all training points to the hyperplane

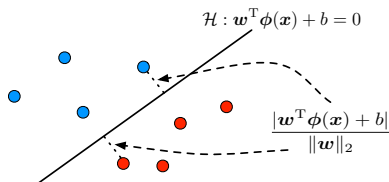
$$\text{MARGIN OF } (\mathbf{w}, b) = \min_n \frac{y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|_2}$$



# Maximizing margin

**Margin:** the *smallest* distance from all training points to the hyperplane

$$\text{MARGIN OF } (\mathbf{w}, b) = \min_n \frac{y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|_2}$$



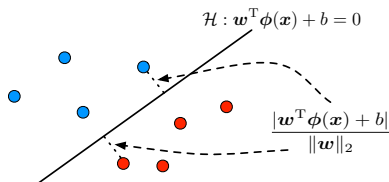
The intuition “**the further away the better**” translates to solving

$$\max_{\mathbf{w}, b} \min_n \frac{y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|}$$

# Maximizing margin

**Margin:** the *smallest* distance from all training points to the hyperplane

$$\text{MARGIN OF } (\mathbf{w}, b) = \min_n \frac{y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|_2}$$



The intuition “**the further away the better**” translates to solving

$$\max_{\mathbf{w}, b} \min_n \frac{y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|} = \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|_2} \min_n y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)$$

# Rescaling

**Note:** rescaling  $(\boldsymbol{w}, b)$  does not change the hyperplane at all.

# Rescaling

**Note:** rescaling  $(\mathbf{w}, b)$  does not change the hyperplane at all.

We can thus always scale  $(\mathbf{w}, b)$  s.t.  $\min_n y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) = 1$

# Rescaling

**Note:** rescaling  $(\mathbf{w}, b)$  does not change the hyperplane at all.

We can thus always scale  $(\mathbf{w}, b)$  s.t.  $\min_n y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) = 1$

The margin then becomes

$$\begin{aligned} & \text{MARGIN OF } (\mathbf{w}, b) \\ &= \frac{1}{\|\mathbf{w}\|_2} \min_n y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \\ &= \frac{1}{\|\mathbf{w}\|_2} \end{aligned}$$

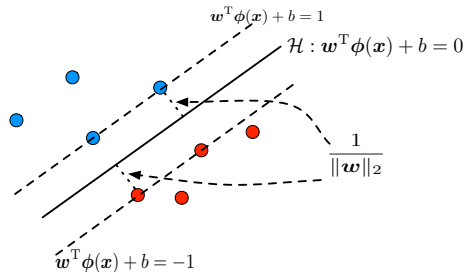
# Rescaling

**Note:** rescaling  $(\mathbf{w}, b)$  does not change the hyperplane at all.

We can thus always scale  $(\mathbf{w}, b)$  s.t.  $\min_n y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) = 1$

The margin then becomes

$$\begin{aligned} \text{MARGIN OF } (\mathbf{w}, b) &= \frac{1}{\|\mathbf{w}\|_2} \min_n y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \\ &= \frac{1}{\|\mathbf{w}\|_2} \end{aligned}$$



# Summary for separable data

For a separable training set, we aim to solve

$$\max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|_2} \quad \text{s.t.} \quad \min_n y_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) = 1$$



## Summary for separable data

For a separable training set, we aim to solve

$$\max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|_2} \quad \text{s.t.} \quad \min_n y_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) = 1$$

This is equivalent to

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1, \quad \forall n \end{aligned}$$

## Summary for separable data

For a separable training set, we aim to solve

$$\max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|_2} \quad \text{s.t.} \quad \min_n y_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) = 1$$

This is equivalent to

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1, \quad \forall n \end{aligned}$$

SVM is thus also called *max-margin* classifier. The constraints above are called *hard-margin* constraints.

## General non-separable case

If data is not linearly separable, the previous constraint

$$y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1, \quad \forall n$$

is obviously *not feasible*.

## General non-separable case

If data is not linearly separable, the previous constraint

$$y_n(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) + b) \geq 1, \quad \forall n$$

is obviously *not feasible*.

To deal with this issue, we relax them to **soft-margin** constraints:

$$y_n(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) + b) \geq 1 - \xi_n, \quad \forall n$$

where we introduce **slack variables**  $\xi_n \geq 0$ .

# SVM Primal formulation

We want  $\xi_n$  to be as small as possible too.

# SVM Primal formulation

We want  $\xi_n$  to be as small as possible too. The objective becomes

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_n\}} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1 - \xi_n, \quad \forall n \\ & \xi_n \geq 0, \quad \forall n \end{aligned}$$

where  $C$  is a hyperparameter to balance the two goals.

# Equivalent form

## Formulation

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_n\}} \quad & C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \leq \xi_n, \quad \forall n \\ & \xi_n \geq 0, \quad \forall n \end{aligned}$$

# Equivalent form

## Formulation

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_n\}} \quad & C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \leq \xi_n, \quad \forall n \\ & \xi_n \geq 0, \quad \forall n \end{aligned}$$

is equivalent to

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_n\}} \quad & C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & \max\{0, 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)\} = \xi_n, \quad \forall n \end{aligned}$$



# Equivalent form

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_n\}} \quad & C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & \max \{0, 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)\} = \xi_n, \quad \forall n \end{aligned}$$

is equivalent to

$$\min_{\mathbf{w}, b} C \sum_n \max \{0, 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)\} + \frac{1}{2} \|\mathbf{w}\|_2^2$$

## Equivalent form

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_n\}} \quad & C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & \max \{0, 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)\} = \xi_n, \quad \forall n \end{aligned}$$

is equivalent to

$$\min_{\mathbf{w}, b} C \sum_n \max \{0, 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)\} + \frac{1}{2} \|\mathbf{w}\|_2^2$$

and

$$\min_{\mathbf{w}, b} \sum_n \max \{0, 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)\} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

with  $\lambda = 1/C$ .

# Equivalent form

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_n\}} \quad & C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & \max \{0, 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)\} = \xi_n, \quad \forall n \end{aligned}$$

is equivalent to

$$\min_{\mathbf{w}, b} C \sum_n \max \{0, 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)\} + \frac{1}{2} \|\mathbf{w}\|_2^2$$

and

$$\min_{\mathbf{w}, b} \sum_n \max \{0, 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)\} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

with  $\lambda = 1/C$ . *This is exactly minimizing L2 regularized hinge loss!*

# Optimization

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_n\}} \quad & C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & 1 - y_n(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) + b) \leq \xi_n, \quad \forall n \\ & \xi_n \geq 0, \quad \forall n \end{aligned}$$

- It is a convex (**quadratic** in fact) problem

# Optimization

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_n\}} \quad & C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \leq \xi_n, \quad \forall n \\ & \xi_n \geq 0, \quad \forall n \end{aligned}$$

- It is a convex (**quadratic** in fact) problem
- thus can apply any convex optimization algorithms, e.g. SGD

# Optimization

$$\begin{aligned}
 \min_{\mathbf{w}, b, \{\xi_n\}} \quad & C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 \\
 \text{s.t.} \quad & 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \leq \xi_n, \quad \forall n \\
 & \xi_n \geq 0, \quad \forall n
 \end{aligned}$$

- It is a convex (**quadratic** in fact) problem
- thus can apply any convex optimization algorithms, e.g. SGD
- there are **more specialized and efficient** algorithms

# Optimization

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_n\}} \quad & C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \leq \xi_n, \quad \forall n \\ & \xi_n \geq 0, \quad \forall n \end{aligned}$$

- It is a convex (**quadratic** in fact) problem
- thus can apply any convex optimization algorithms, e.g. SGD
- there are **more specialized and efficient** algorithms
- but usually we apply kernel trick, which requires solving the *dual problem*

# Outline

- 1 Review of last lecture
- 2 Support vector machines (primal formulation)
- 3 A detour of Lagrangian duality**
- 4 Support vector machines (dual formulation)



# Lagrangian duality

Extremely important and powerful tool in analyzing optimizations

# Lagrangian duality

Extremely important and powerful tool in analyzing optimizations

We will introduce basic concepts and derive the **KKT conditions**

# Lagrangian duality

Extremely important and powerful tool in analyzing optimizations

We will introduce basic concepts and derive the **KKT conditions**

Applying it to SVM reveals an important aspect of the algorithm

# Primal problem

Suppose we want to solve

$$\min_{\mathbf{w}} F(\mathbf{w}) \quad \text{s.t.} \quad h_j(\mathbf{w}) \leq 0 \quad \forall j \in [J]$$

where functions  $h_1, \dots, h_J$  define  $J$  **constraints**.

## Primal problem

Suppose we want to solve

$$\min_{\mathbf{w}} F(\mathbf{w}) \quad \text{s.t.} \quad h_j(\mathbf{w}) \leq 0 \quad \forall j \in [J]$$

where functions  $h_1, \dots, h_J$  define  $J$  **constraints**.

SVM primal formulation is clearly of this form with  $J = 2N$  constraints:

$$F(\mathbf{w}, b, \{\xi_n\}) = C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$h_n(\mathbf{w}, b, \{\xi_n\}) = 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \xi_n \quad \forall n \in [N]$$

$$h_{N+n}(\mathbf{w}, b, \{\xi_n\}) = -\xi_n \quad \forall n \in [N]$$

# Lagrangian

The **Lagrangian** of the previous problem is defined as:

$$L(\mathbf{w}, \{\lambda_j\}) = F(\mathbf{w}) + \sum_{j=1}^J \lambda_j h_j(\mathbf{w})$$

where  $\lambda_1, \dots, \lambda_J \geq 0$  are called **Lagrangian multipliers**.

# Lagrangian

The **Lagrangian** of the previous problem is defined as:

$$L(\mathbf{w}, \{\lambda_j\}) = F(\mathbf{w}) + \sum_{j=1}^J \lambda_j h_j(\mathbf{w})$$

where  $\lambda_1, \dots, \lambda_J \geq 0$  are called **Lagrangian multipliers**.

Note that

$$\max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\}) = \begin{cases} F(\mathbf{w}) & \text{if } h_j(\mathbf{w}) \leq 0 \quad \forall j \in [J] \\ +\infty & \text{else} \end{cases}$$

# Lagrangian

The **Lagrangian** of the previous problem is defined as:

$$L(\mathbf{w}, \{\lambda_j\}) = F(\mathbf{w}) + \sum_{j=1}^J \lambda_j h_j(\mathbf{w})$$

where  $\lambda_1, \dots, \lambda_J \geq 0$  are called **Lagrangian multipliers**.

Note that

$$\max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\}) = \begin{cases} F(\mathbf{w}) & \text{if } h_j(\mathbf{w}) \leq 0 \quad \forall j \in [J] \\ +\infty & \text{else} \end{cases}$$

and thus,

$$\min_{\mathbf{w}} \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\}) \iff \min_{\mathbf{w}} F(\mathbf{w}) \quad \text{s.t.} \quad h_j(\mathbf{w}) \leq 0 \quad \forall j \in [J]$$



# Duality

We define the **dual problem** by swapping the min and max:

$$\max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\})$$

# Duality

We define the **dual problem** by swapping the min and max:

$$\max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\})$$

*How are the primal and dual connected?*

# Duality

We define the **dual problem** by swapping the min and max:

$$\max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\})$$

*How are the primal and dual connected?* Let  $\mathbf{w}^*$  and  $\{\lambda_j\}^*$  be the primal and dual solutions respectively

# Duality

We define the **dual problem** by swapping the min and max:

$$\max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\})$$

*How are the primal and dual connected?* Let  $\mathbf{w}^*$  and  $\{\lambda_j\}_j^*$  be the primal and dual solutions respectively, then

$$\max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\}) = \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j^*\})$$

# Duality

We define the **dual problem** by swapping the min and max:

$$\max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\})$$

*How are the primal and dual connected?* Let  $\mathbf{w}^*$  and  $\{\lambda_j\}^*$  be the primal and dual solutions respectively, then

$$\max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\}) = \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j^*\}) \leq L(\mathbf{w}^*, \{\lambda_j^*\})$$

# Duality

We define the **dual problem** by swapping the min and max:

$$\max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\})$$

*How are the primal and dual connected?* Let  $\mathbf{w}^*$  and  $\{\lambda_j\}_j^*$  be the primal and dual solutions respectively, then

$$\begin{aligned} \max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\}) &= \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j^*\}) \leq L(\mathbf{w}^*, \{\lambda_j^*\}) \\ &\leq \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}^*, \{\lambda_j\}) \end{aligned}$$

# Duality

We define the **dual problem** by swapping the min and max:

$$\max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\})$$

*How are the primal and dual connected?* Let  $\mathbf{w}^*$  and  $\{\lambda_j\}_j^*$  be the primal and dual solutions respectively, then

$$\begin{aligned} \max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\}) &= \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j^*\}) \leq L(\mathbf{w}^*, \{\lambda_j^*\}) \\ &\leq \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}^*, \{\lambda_j\}) = \min_{\mathbf{w}} \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\}) \end{aligned}$$

# Duality

We define the **dual problem** by swapping the min and max:

$$\max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\})$$

*How are the primal and dual connected?* Let  $\mathbf{w}^*$  and  $\{\lambda_j\}^*$  be the primal and dual solutions respectively, then

$$\begin{aligned} \max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\}) &= \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j^*\}) \leq L(\mathbf{w}^*, \{\lambda_j^*\}) \\ &\leq \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}^*, \{\lambda_j\}) = \min_{\mathbf{w}} \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\}) \end{aligned}$$

This is called “**weak duality**”.



# Strong duality

When  $F, h_1, \dots, h_m$  are convex, under some mild conditions:

$$\min_{\mathbf{w}} \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\}) = \max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\})$$

# Strong duality

When  $F, h_1, \dots, h_m$  are convex, under some mild conditions:

$$\min_{\mathbf{w}} \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\}) = \max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\})$$

This is called “**strong duality**”.

# Deriving the Karush-Kuhn-Tucker (KKT) conditions

**Observe that if strong duality holds:**

$$F(\mathbf{w}^*) = \min_{\mathbf{w}} \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\}) = \max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\})$$

# Deriving the Karush-Kuhn-Tucker (KKT) conditions

**Observe that if strong duality holds:**

$$\begin{aligned} F(\mathbf{w}^*) &= \min_{\mathbf{w}} \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\}) = \max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\}) \\ &= \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j^*\}) \end{aligned}$$

# Deriving the Karush-Kuhn-Tucker (KKT) conditions

Observe that if strong duality holds:

$$\begin{aligned} F(\mathbf{w}^*) &= \min_{\mathbf{w}} \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\}) = \max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\}) \\ &= \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j^*\}) \leq L(\mathbf{w}^*, \{\lambda_j^*\}) \end{aligned}$$

# Deriving the Karush-Kuhn-Tucker (KKT) conditions

Observe that if strong duality holds:

$$\begin{aligned}
 F(\mathbf{w}^*) &= \min_{\mathbf{w}} \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\}) = \max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\}) \\
 &= \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j^*\}) \leq L(\mathbf{w}^*, \{\lambda_j^*\}) = F(\mathbf{w}^*) + \sum_{j=1}^J \lambda_j^* h_j(\mathbf{w}^*)
 \end{aligned}$$

# Deriving the Karush-Kuhn-Tucker (KKT) conditions

Observe that if strong duality holds:

$$\begin{aligned}
 F(\mathbf{w}^*) &= \min_{\mathbf{w}} \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\}) = \max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\}) \\
 &= \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j^*\}) \leq L(\mathbf{w}^*, \{\lambda_j^*\}) = F(\mathbf{w}^*) + \sum_{j=1}^J \lambda_j^* h_j(\mathbf{w}^*) \leq F(\mathbf{w}^*)
 \end{aligned}$$

# Deriving the Karush-Kuhn-Tucker (KKT) conditions

Observe that if strong duality holds:

$$\begin{aligned}
 F(\mathbf{w}^*) &= \min_{\mathbf{w}} \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\}) = \max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\}) \\
 &= \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j^*\}) \leq L(\mathbf{w}^*, \{\lambda_j^*\}) = F(\mathbf{w}^*) + \sum_{j=1}^J \lambda_j^* h_j(\mathbf{w}^*) \leq F(\mathbf{w}^*)
 \end{aligned}$$

Implications:

- *all inequalities above have to be equalities!*



# Deriving the Karush-Kuhn-Tucker (KKT) conditions

Observe that if strong duality holds:

$$\begin{aligned}
 F(\mathbf{w}^*) &= \min_{\mathbf{w}} \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\}) = \max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\}) \\
 &= \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j^*\}) \leq L(\mathbf{w}^*, \{\lambda_j^*\}) = F(\mathbf{w}^*) + \sum_{j=1}^J \lambda_j^* h_j(\mathbf{w}^*) \leq F(\mathbf{w}^*)
 \end{aligned}$$

Implications:

- *all inequalities above have to be equalities!*
- last equality implies  $\lambda_j^* h_j(\mathbf{w}^*) = 0$  for all  $j \in [J]$

# Deriving the Karush-Kuhn-Tucker (KKT) conditions

Observe that if strong duality holds:

$$\begin{aligned}
 F(\mathbf{w}^*) &= \min_{\mathbf{w}} \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\}) = \max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\}) \\
 &= \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j^*\}) \leq L(\mathbf{w}^*, \{\lambda_j^*\}) = F(\mathbf{w}^*) + \sum_{j=1}^J \lambda_j^* h_j(\mathbf{w}^*) \leq F(\mathbf{w}^*)
 \end{aligned}$$

Implications:

- *all inequalities above have to be equalities!*
- last equality implies  $\lambda_j^* h_j(\mathbf{w}^*) = 0$  for all  $j \in [J]$
- equality  $\min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j^*\}) = L(\mathbf{w}^*, \{\lambda_j^*\})$  implies  $\mathbf{w}^*$  is a **minimizer** of  $L(\mathbf{w}, \{\lambda_j^*\})$

# Deriving the Karush-Kuhn-Tucker (KKT) conditions

Observe that if strong duality holds:

$$\begin{aligned}
 F(\mathbf{w}^*) &= \min_{\mathbf{w}} \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\}) = \max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\}) \\
 &= \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j^*\}) \leq L(\mathbf{w}^*, \{\lambda_j^*\}) = F(\mathbf{w}^*) + \sum_{j=1}^J \lambda_j^* h_j(\mathbf{w}^*) \leq F(\mathbf{w}^*)
 \end{aligned}$$

Implications:

- *all inequalities above have to be equalities!*
- last equality implies  $\lambda_j^* h_j(\mathbf{w}^*) = 0$  for all  $j \in [J]$
- equality  $\min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j^*\}) = L(\mathbf{w}^*, \{\lambda_j^*\})$  implies  $\mathbf{w}^*$  is a **minimizer** of  $L(\mathbf{w}, \{\lambda_j^*\})$  and thus has **zero gradient**:

$$\nabla_{\mathbf{w}} L(\mathbf{w}^*, \{\lambda_j^*\}) = \nabla F(\mathbf{w}^*) + \sum_{j=1}^J \lambda_j^* \nabla h_j(\mathbf{w}^*) = \mathbf{0}$$

# The Karush-Kuhn-Tucker (KKT) conditions

If  $w^*$  and  $\{\lambda_j^*\}$  are the primal and dual solution respectively, then:

# The Karush-Kuhn-Tucker (KKT) conditions

If  $\mathbf{w}^*$  and  $\{\lambda_j^*\}$  are the primal and dual solution respectively, then:

**Stationarity:**

$$\nabla_{\mathbf{w}} L(\mathbf{w}^*, \{\lambda_j^*\}) = \nabla F(\mathbf{w}^*) + \sum_{j=1}^J \lambda_j^* \nabla h_j(\mathbf{w}^*) = \mathbf{0}$$

# The Karush-Kuhn-Tucker (KKT) conditions

If  $\mathbf{w}^*$  and  $\{\lambda_j^*\}$  are the primal and dual solution respectively, then:

## Stationarity:

$$\nabla_{\mathbf{w}} L(\mathbf{w}^*, \{\lambda_j^*\}) = \nabla F(\mathbf{w}^*) + \sum_{j=1}^J \lambda_j^* \nabla h_j(\mathbf{w}^*) = \mathbf{0}$$

## Complementary slackness:

$$\lambda_j^* h_j(\mathbf{w}^*) = 0 \quad \text{for all } j \in [J]$$

# The Karush-Kuhn-Tucker (KKT) conditions

If  $\mathbf{w}^*$  and  $\{\lambda_j^*\}$  are the primal and dual solution respectively, then:

## Stationarity:

$$\nabla_{\mathbf{w}} L(\mathbf{w}^*, \{\lambda_j^*\}) = \nabla F(\mathbf{w}^*) + \sum_{j=1}^J \lambda_j^* \nabla h_j(\mathbf{w}^*) = \mathbf{0}$$

## Complementary slackness:

$$\lambda_j^* h_j(\mathbf{w}^*) = 0 \quad \text{for all } j \in [J]$$

## Feasibility:

$$h_j(\mathbf{w}^*) \leq 0 \quad \text{and} \quad \lambda_j^* \geq 0 \quad \text{for all } j \in [J]$$

# The Karush-Kuhn-Tucker (KKT) conditions

If  $\mathbf{w}^*$  and  $\{\lambda_j^*\}$  are the primal and dual solution respectively, then:

## Stationarity:

$$\nabla_{\mathbf{w}} L(\mathbf{w}^*, \{\lambda_j^*\}) = \nabla F(\mathbf{w}^*) + \sum_{j=1}^J \lambda_j^* \nabla h_j(\mathbf{w}^*) = \mathbf{0}$$

## Complementary slackness:

$$\lambda_j^* h_j(\mathbf{w}^*) = 0 \quad \text{for all } j \in [J]$$

## Feasibility:

$$h_j(\mathbf{w}^*) \leq 0 \quad \text{and} \quad \lambda_j^* \geq 0 \quad \text{for all } j \in [J]$$

These are *necessary conditions*.



# The Karush-Kuhn-Tucker (KKT) conditions

If  $\mathbf{w}^*$  and  $\{\lambda_j^*\}$  are the primal and dual solution respectively, then:

## Stationarity:

$$\nabla_{\mathbf{w}} L(\mathbf{w}^*, \{\lambda_j^*\}) = \nabla F(\mathbf{w}^*) + \sum_{j=1}^J \lambda_j^* \nabla h_j(\mathbf{w}^*) = \mathbf{0}$$

## Complementary slackness:

$$\lambda_j^* h_j(\mathbf{w}^*) = 0 \quad \text{for all } j \in [J]$$

## Feasibility:

$$h_j(\mathbf{w}^*) \leq 0 \quad \text{and} \quad \lambda_j^* \geq 0 \quad \text{for all } j \in [J]$$

These are *necessary conditions*. They are also *sufficient* when  $F$  is convex and  $h_1, \dots, h_J$  are continuously differentiable convex functions.

# Outline

- 1 Review of last lecture
- 2 Support vector machines (primal formulation)
- 3 A detour of Lagrangian duality
- 4 Support vector machines (dual formulation)**

# Writing down the Lagrangian

Recall the primal formulation

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_n\}} \quad & C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & 1 - y_n [\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) + b] \leq \xi_n, \quad \forall n \\ & \xi_n \geq 0, \quad \forall n \end{aligned}$$

# Writing down the Lagrangian

Recall the primal formulation

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_n\}} \quad & C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & 1 - y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \leq \xi_n, \quad \forall n \\ & \xi_n \geq 0, \quad \forall n \end{aligned}$$

**Lagrangian** is

$$\begin{aligned} L(\mathbf{w}, b, \{\xi_n\}, \{\alpha_n\}, \{\lambda_n\}) = & C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_n \lambda_n \xi_n \\ & + \sum_n \alpha_n (1 - y_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \xi_n) \end{aligned}$$

where  $\alpha_1, \dots, \alpha_N \geq 0$  and  $\lambda_1, \dots, \lambda_N \geq 0$  are Lagrangian multipliers.

## Applying the stationarity condition

$$L = C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_n \lambda_n \xi_n + \sum_n \alpha_n (1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \xi_n)$$

$\exists$  primal and dual variables  $\mathbf{w}, b, \{\xi_n\}, \{\alpha_n\}, \{\lambda_n\}$  s.t.  $\nabla_{\mathbf{w}, b, \{\xi_n\}} L = \mathbf{0}$ ,

## Applying the stationarity condition

$$L = C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_n \lambda_n \xi_n + \sum_n \alpha_n (1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \xi_n)$$

$\exists$  primal and dual variables  $\mathbf{w}, b, \{\xi_n\}, \{\alpha_n\}, \{\lambda_n\}$  s.t.  $\nabla_{\mathbf{w}, b, \{\xi_n\}} L = \mathbf{0}$ , which means

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_n y_n \alpha_n \phi(\mathbf{x}_n) = \mathbf{0}$$

## Applying the stationarity condition

$$L = C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_n \lambda_n \xi_n + \sum_n \alpha_n (1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \xi_n)$$

$\exists$  primal and dual variables  $\mathbf{w}, b, \{\xi_n\}, \{\alpha_n\}, \{\lambda_n\}$  s.t.  $\nabla_{\mathbf{w}, b, \{\xi_n\}} L = \mathbf{0}$ , which means

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_n y_n \alpha_n \phi(\mathbf{x}_n) = \mathbf{0} \quad \implies \quad \mathbf{w} = \sum_n y_n \alpha_n \phi(\mathbf{x}_n)$$

## Applying the stationarity condition

$$L = C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_n \lambda_n \xi_n + \sum_n \alpha_n (1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \xi_n)$$

$\exists$  primal and dual variables  $\mathbf{w}, b, \{\xi_n\}, \{\alpha_n\}, \{\lambda_n\}$  s.t.  $\nabla_{\mathbf{w}, b, \{\xi_n\}} L = \mathbf{0}$ , which means

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_n y_n \alpha_n \phi(\mathbf{x}_n) = \mathbf{0} \quad \Longrightarrow \quad \mathbf{w} = \sum_n y_n \alpha_n \phi(\mathbf{x}_n)$$

$$\frac{\partial L}{\partial b} = - \sum_n \alpha_n y_n = 0$$



## Applying the stationarity condition

$$L = C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_n \lambda_n \xi_n + \sum_n \alpha_n (1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \xi_n)$$

$\exists$  primal and dual variables  $\mathbf{w}, b, \{\xi_n\}, \{\alpha_n\}, \{\lambda_n\}$  s.t.  $\nabla_{\mathbf{w}, b, \{\xi_n\}} L = \mathbf{0}$ , which means

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_n y_n \alpha_n \phi(\mathbf{x}_n) = \mathbf{0} \quad \implies \quad \mathbf{w} = \sum_n y_n \alpha_n \phi(\mathbf{x}_n)$$

$$\frac{\partial L}{\partial b} = - \sum_n \alpha_n y_n = 0 \quad \text{and} \quad \frac{\partial L}{\partial \xi_n} = C - \lambda_n - \alpha_n = 0, \quad \forall n$$

## Rewrite the Lagrangian in terms of dual variables

Replacing  $w$  by  $\sum_n y_n \alpha_n \phi(\mathbf{x}_n)$  in the Lagrangian gives

$$L = C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_n \lambda_n \xi_n + \sum_n \alpha_n (1 - y_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \xi_n)$$

# Rewrite the Lagrangian in terms of dual variables

Replacing  $w$  by  $\sum_n y_n \alpha_n \phi(\mathbf{x}_n)$  in the Lagrangian gives

$$\begin{aligned}
 L &= C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_n \lambda_n \xi_n + \sum_n \alpha_n (1 - y_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \xi_n) \\
 &= C \sum_n \xi_n + \frac{1}{2} \left\| \sum_n y_n \alpha_n \phi(\mathbf{x}_n) \right\|_2^2 - \sum_n \lambda_n \xi_n + \\
 &\quad \sum_n \alpha_n \left( 1 - y_n \left( \left( \sum_m y_m \alpha_m \phi(\mathbf{x}_m) \right)^T \phi(\mathbf{x}_n) + b \right) - \xi_n \right)
 \end{aligned}$$

# Rewrite the Lagrangian in terms of dual variables

Replacing  $w$  by  $\sum_n y_n \alpha_n \phi(\mathbf{x}_n)$  in the Lagrangian gives

$$\begin{aligned}
 L &= C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_n \lambda_n \xi_n + \sum_n \alpha_n (1 - y_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \xi_n) \\
 &= C \sum_n \xi_n + \frac{1}{2} \left\| \sum_n y_n \alpha_n \phi(\mathbf{x}_n) \right\|_2^2 - \sum_n \lambda_n \xi_n + \\
 &\quad \sum_n \alpha_n \left( 1 - y_n \left( \left( \sum_m y_m \alpha_m \phi(\mathbf{x}_m) \right)^T \phi(\mathbf{x}_n) + b \right) - \xi_n \right) \\
 &= \sum_n \alpha_n + \frac{1}{2} \left\| \sum_n y_n \alpha_n \phi(\mathbf{x}_n) \right\|_2^2 - \sum_{m,n} \alpha_n \alpha_m y_m y_n \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n) \\
 &\quad \left( \sum_n \alpha_n y_n = 0 \text{ and } C = \lambda_n + \alpha_n \right)
 \end{aligned}$$

# Rewrite the Lagrangian in terms of dual variables

Replacing  $w$  by  $\sum_n y_n \alpha_n \phi(\mathbf{x}_n)$  in the Lagrangian gives

$$\begin{aligned}
 L &= C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_n \lambda_n \xi_n + \sum_n \alpha_n (1 - y_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \xi_n) \\
 &= C \sum_n \xi_n + \frac{1}{2} \left\| \sum_n y_n \alpha_n \phi(\mathbf{x}_n) \right\|_2^2 - \sum_n \lambda_n \xi_n + \\
 &\quad \sum_n \alpha_n \left( 1 - y_n \left( \left( \sum_m y_m \alpha_m \phi(\mathbf{x}_m) \right)^T \phi(\mathbf{x}_n) + b \right) - \xi_n \right) \\
 &= \sum_n \alpha_n + \frac{1}{2} \left\| \sum_n y_n \alpha_n \phi(\mathbf{x}_n) \right\|_2^2 - \sum_{m,n} \alpha_n \alpha_m y_m y_n \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n) \\
 &\qquad\qquad\qquad (\sum_n \alpha_n y_n = 0 \text{ and } C = \lambda_n + \alpha_n) \\
 &= \sum_n \alpha_n - \frac{1}{2} \sum_{m,n} \alpha_n \alpha_m y_m y_n \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)
 \end{aligned}$$

# The dual formulation

To find the dual solutions, it amounts to solving

$$\begin{aligned} \max_{\{\alpha_n\}, \{\lambda_n\}} \quad & \sum_n \alpha_n - \frac{1}{2} \sum_{m,n} y_m y_n \alpha_m \alpha_n \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n) \\ \text{s.t.} \quad & \sum_n \alpha_n y_n = 0 \\ & C - \lambda_n - \alpha_n = 0, \quad \alpha_n \geq 0, \quad \lambda_n \geq 0, \quad \forall n \end{aligned}$$

# The dual formulation

To find the dual solutions, it amounts to solving

$$\begin{aligned} \max_{\{\alpha_n\}, \{\lambda_n\}} \quad & \sum_n \alpha_n - \frac{1}{2} \sum_{m,n} y_m y_n \alpha_m \alpha_n \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n) \\ \text{s.t.} \quad & \sum_n \alpha_n y_n = 0 \\ & C - \lambda_n - \alpha_n = 0, \quad \alpha_n \geq 0, \quad \lambda_n \geq 0, \quad \forall n \end{aligned}$$

Note the last three constraints can be written as  $0 \leq \alpha_n \leq C$  for all  $n$ .

# The dual formulation

To find the dual solutions, it amounts to solving

$$\begin{aligned} \max_{\{\alpha_n\}, \{\lambda_n\}} \quad & \sum_n \alpha_n - \frac{1}{2} \sum_{m,n} y_m y_n \alpha_m \alpha_n \phi(\mathbf{x}_m)^\top \phi(\mathbf{x}_n) \\ \text{s.t.} \quad & \sum_n \alpha_n y_n = 0 \\ & C - \lambda_n - \alpha_n = 0, \quad \alpha_n \geq 0, \quad \lambda_n \geq 0, \quad \forall n \end{aligned}$$

Note the last three constraints can be written as  $0 \leq \alpha_n \leq C$  for all  $n$ . So the final **dual formulation of SVM** is:

$$\begin{aligned} \max_{\{\alpha_n\}} \quad & \sum_n \alpha_n - \frac{1}{2} \sum_{m,n} y_m y_n \alpha_m \alpha_n \phi(\mathbf{x}_m)^\top \phi(\mathbf{x}_n) \\ \text{s.t.} \quad & \sum_n \alpha_n y_n = 0 \quad \text{and} \quad 0 \leq \alpha_n \leq C, \quad \forall n \end{aligned}$$



# Kernelizing SVM

Now it is clear that with a **kernel function**  $k$  for the mapping  $\phi$ , we can kernelize SVM as:

$$\begin{aligned} \max_{\{\alpha_n\}} \quad & \sum_n \alpha_n - \frac{1}{2} \sum_{m,n} y_m y_n \alpha_m \alpha_n k(\mathbf{x}_m, \mathbf{x}_n) \\ \text{s.t.} \quad & \sum_n \alpha_n y_n = 0 \quad \text{and} \quad 0 \leq \alpha_n \leq C, \quad \forall n \end{aligned}$$

Again, no need to compute  $\phi(\mathbf{x})$ .

# Kernelizing SVM

Now it is clear that with a **kernel function**  $k$  for the mapping  $\phi$ , we can kernelize SVM as:

$$\begin{aligned} \max_{\{\alpha_n\}} \quad & \sum_n \alpha_n - \frac{1}{2} \sum_{m,n} y_m y_n \alpha_m \alpha_n k(\mathbf{x}_m, \mathbf{x}_n) \\ \text{s.t.} \quad & \sum_n \alpha_n y_n = 0 \quad \text{and} \quad 0 \leq \alpha_n \leq C, \quad \forall n \end{aligned}$$

Again, no need to compute  $\phi(\mathbf{x})$ . It is a **quadratic program** and many efficient optimization algorithms exist.

# Recover the primal solution

But how do we predict given the dual solution  $\{\alpha_n^*\}$ ?

# Recover the primal solution

But how do we predict given the dual solution  $\{\alpha_n^*\}$ ? Need to figure out the primal solution  $w^*$  and  $b^*$ .

# Recover the primal solution

But how do we predict given the dual solution  $\{\alpha_n^*\}$ ? Need to figure out the primal solution  $\mathbf{w}^*$  and  $b^*$ .

Based on previous observation,

$$\mathbf{w}^* = \sum_n \alpha_n^* y_n \phi(\mathbf{x}_n)$$

# Recover the primal solution

But how do we predict given the dual solution  $\{\alpha_n^*\}$ ? Need to figure out the primal solution  $\mathbf{w}^*$  and  $b^*$ .

Based on previous observation,

$$\mathbf{w}^* = \sum_n \alpha_n^* y_n \phi(\mathbf{x}_n) = \sum_{n:\alpha_n > 0} \alpha_n^* y_n \phi(\mathbf{x}_n)$$

## Recover the primal solution

But how do we predict given the dual solution  $\{\alpha_n^*\}$ ? Need to figure out the primal solution  $\mathbf{w}^*$  and  $b^*$ .

Based on previous observation,

$$\mathbf{w}^* = \sum_n \alpha_n^* y_n \phi(\mathbf{x}_n) = \sum_{n:\alpha_n^* > 0} \alpha_n^* y_n \phi(\mathbf{x}_n)$$

A point with  $\alpha_n^* > 0$  is called a “**support vector**”. Hence the name SVM.

## Recover the primal solution

But how do we predict given the dual solution  $\{\alpha_n^*\}$ ? Need to figure out the primal solution  $\mathbf{w}^*$  and  $b^*$ .

Based on previous observation,

$$\mathbf{w}^* = \sum_n \alpha_n^* y_n \phi(\mathbf{x}_n) = \sum_{n:\alpha_n^* > 0} \alpha_n^* y_n \phi(\mathbf{x}_n)$$

A point with  $\alpha_n^* > 0$  is called a “**support vector**”. Hence the name SVM.

To identify  $b$ , we need to apply complementary slackness.



## Applying complementary slackness

For all  $n$  we should have

$$\lambda_n^* \xi_n^* = 0, \quad \alpha_n^* \left( 1 - \xi_n^* - y_n (\mathbf{w}^{*\top} \boldsymbol{\phi}(\mathbf{x}_n) + b^*) \right) = 0$$

## Applying complementary slackness

For all  $n$  we should have

$$\lambda_n^* \xi_n^* = 0, \quad \alpha_n^* \left( 1 - \xi_n^* - y_n (\mathbf{w}^{*\top} \phi(\mathbf{x}_n) + b^*) \right) = 0$$

For some support vector  $\phi(\mathbf{x}_n)$  if we have  $0 < \alpha_n^* < C$ , then

$$\lambda_n^* = C - \alpha_n^* > 0$$

## Applying complementary slackness

For all  $n$  we should have

$$\lambda_n^* \xi_n^* = 0, \quad \alpha_n^* \left( 1 - \xi_n^* - y_n (\mathbf{w}^{*\top} \phi(\mathbf{x}_n) + b^*) \right) = 0$$

For some support vector  $\phi(\mathbf{x}_n)$  if we have  $0 < \alpha_n^* < C$ , then

$$\lambda_n^* = C - \alpha_n^* > 0$$

With the first condition we know  $\xi_n^* = 0$ .

## Applying complementary slackness

For all  $n$  we should have

$$\lambda_n^* \xi_n^* = 0, \quad \alpha_n^* \left( 1 - \xi_n^* - y_n (\mathbf{w}^{*\top} \phi(\mathbf{x}_n) + b^*) \right) = 0$$

For some support vector  $\phi(\mathbf{x}_n)$  if we have  $0 < \alpha_n^* < C$ , then

$$\lambda_n^* = C - \alpha_n^* > 0$$

With the first condition we know  $\xi_n^* = 0$ .

With the second condition we know  $1 = y_n (\mathbf{w}^{*\top} \phi(\mathbf{x}_n) + b^*)$

## Applying complementary slackness

For all  $n$  we should have

$$\lambda_n^* \xi_n^* = 0, \quad \alpha_n^* \left( 1 - \xi_n^* - y_n (\mathbf{w}^{*\top} \phi(\mathbf{x}_n) + b^*) \right) = 0$$

For some support vector  $\phi(\mathbf{x}_n)$  if we have  $0 < \alpha_n^* < C$ , then

$$\lambda_n^* = C - \alpha_n^* > 0$$

With the first condition we know  $\xi_n^* = 0$ .

With the second condition we know  $1 = y_n (\mathbf{w}^{*\top} \phi(\mathbf{x}_n) + b^*)$  and thus

$$b^* = y_n - \mathbf{w}^{*\top} \phi(\mathbf{x}_n)$$

## Applying complementary slackness

For all  $n$  we should have

$$\lambda_n^* \xi_n^* = 0, \quad \alpha_n^* \left( 1 - \xi_n^* - y_n (\mathbf{w}^{*\top} \phi(\mathbf{x}_n) + b^*) \right) = 0$$

For some support vector  $\phi(\mathbf{x}_n)$  if we have  $0 < \alpha_n^* < C$ , then

$$\lambda_n^* = C - \alpha_n^* > 0$$

With the first condition we know  $\xi_n^* = 0$ .

With the second condition we know  $1 = y_n (\mathbf{w}^{*\top} \phi(\mathbf{x}_n) + b^*)$  and thus

$$b^* = y_n - \mathbf{w}^{*\top} \phi(\mathbf{x}_n) = y_n - \sum_m y_m \alpha_m^* k(\mathbf{x}_m, \mathbf{x}_n)$$

## Applying complementary slackness

For all  $n$  we should have

$$\lambda_n^* \xi_n^* = 0, \quad \alpha_n^* \left( 1 - \xi_n^* - y_n (\mathbf{w}^{*\top} \phi(\mathbf{x}_n) + b^*) \right) = 0$$

For some support vector  $\phi(\mathbf{x}_n)$  if we have  $0 < \alpha_n^* < C$ , then

$$\lambda_n^* = C - \alpha_n^* > 0$$

With the first condition we know  $\xi_n^* = 0$ .

With the second condition we know  $1 = y_n (\mathbf{w}^{*\top} \phi(\mathbf{x}_n) + b^*)$  and thus

$$b^* = y_n - \mathbf{w}^{*\top} \phi(\mathbf{x}_n) = y_n - \sum_m y_m \alpha_m^* k(\mathbf{x}_m, \mathbf{x}_n)$$

The prediction on a new point  $\mathbf{x}$  is therefore

$$\text{SGN} \left( \mathbf{w}^{*\top} \phi(\mathbf{x}) + b^* \right) = \text{SGN} \left( \sum_m y_m \alpha_m^* k(\mathbf{x}_m, \mathbf{x}) + b^* \right)$$

# Geometric interpretation of support vectors

A support vector satisfies  $\alpha_n^* \neq 0$  and

$$1 - \xi_n^* - y_n(\mathbf{w}^{*\top} \phi(\mathbf{x}_n) + b^*) = 0$$



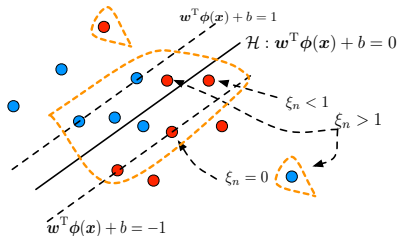
# Geometric interpretation of support vectors

A support vector satisfies  $\alpha_n^* \neq 0$  and

$$1 - \xi_n^* - y_n(\mathbf{w}^{*\top} \phi(\mathbf{x}_n) + b^*) = 0$$

When

- $\xi_n^* = 0$ ,  $y_n(\mathbf{w}^{*\top} \phi(\mathbf{x}_n) + b^*) = 1$   
and thus the point is  $1/\|\mathbf{w}^*\|_2$   
away from the hyperplane.



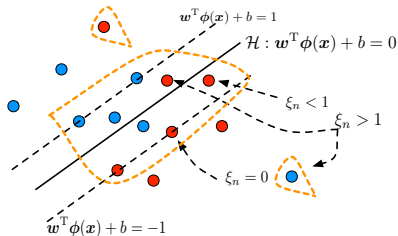
# Geometric interpretation of support vectors

A support vector satisfies  $\alpha_n^* \neq 0$  and

$$1 - \xi_n^* - y_n(\mathbf{w}^{*\top} \boldsymbol{\phi}(\mathbf{x}_n) + b^*) = 0$$

When

- $\xi_n^* = 0$ ,  $y_n(\mathbf{w}^{*\top} \boldsymbol{\phi}(\mathbf{x}_n) + b^*) = 1$  and thus the point is  $1/\|\mathbf{w}^*\|_2$  away from the hyperplane.
- $\xi_n^* < 1$ , the point is classified correctly but does not satisfy the large margin constraint.



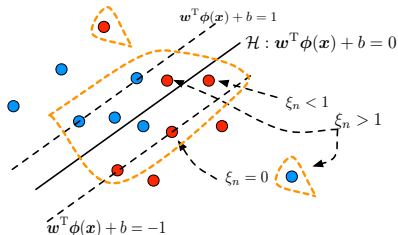
# Geometric interpretation of support vectors

A support vector satisfies  $\alpha_n^* \neq 0$  and

$$1 - \xi_n^* - y_n(\mathbf{w}^{*\top} \phi(\mathbf{x}_n) + b^*) = 0$$

When

- $\xi_n^* = 0$ ,  $y_n(\mathbf{w}^{*\top} \phi(\mathbf{x}_n) + b^*) = 1$  and thus the point is  $1/\|\mathbf{w}^*\|_2$  away from the hyperplane.
- $\xi_n^* < 1$ , the point is classified correctly but does not satisfy the large margin constraint.
- $\xi_n^* > 1$ , the point is misclassified.



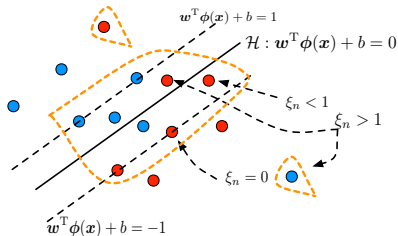
# Geometric interpretation of support vectors

A support vector satisfies  $\alpha_n^* \neq 0$  and

$$1 - \xi_n^* - y_n(\mathbf{w}^{*\top} \boldsymbol{\phi}(\mathbf{x}_n) + b^*) = 0$$

When

- $\xi_n^* = 0$ ,  $y_n(\mathbf{w}^{*\top} \boldsymbol{\phi}(\mathbf{x}_n) + b^*) = 1$  and thus the point is  $1/\|\mathbf{w}^*\|_2$  away from the hyperplane.
- $\xi_n^* < 1$ , the point is classified correctly but does not satisfy the large margin constraint.
- $\xi_n^* > 1$ , the point is misclassified.



Support vectors (circled with the orange line) are *the only points that matter!*

## An example

One drawback of kernel method: **non-parametric**, need to keep all training points potentially

## An example

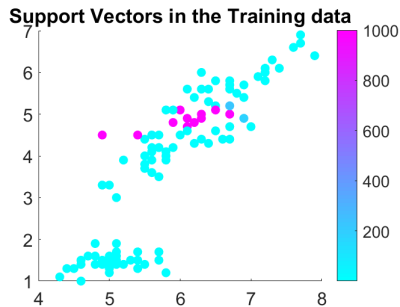
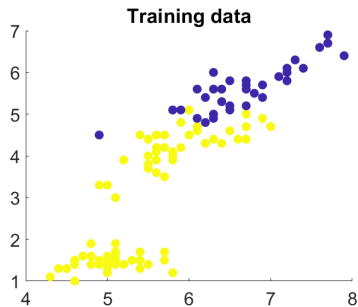
One drawback of kernel method: **non-parametric**, need to keep all training points potentially

For SVM, very often **#support vectors**  $\ll N$

## An example

One drawback of kernel method: **non-parametric**, need to keep all training points potentially

For SVM, very often  $\# \text{support vectors} \ll N$



# Summary

SVM: **max-margin linear classifier**



# Summary

SVM: **max-margin linear classifier**

**Primal** (equivalent to minimizing L2 regularized hinge loss):

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_n\}} \quad & C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \leq \xi_n, \quad \forall n \\ & \xi_n \geq 0, \quad \forall n \end{aligned}$$

# Summary

SVM: **max-margin linear classifier**

**Primal** (equivalent to minimizing L2 regularized hinge loss):

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_n\}} \quad & C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & 1 - y_n(\mathbf{w}^\top \phi(\mathbf{x}_n) + b) \leq \xi_n, \quad \forall n \\ & \xi_n \geq 0, \quad \forall n \end{aligned}$$

**Dual** (kernelizable, reveals what training points are support vectors):

$$\begin{aligned} \max_{\{\alpha_n\}} \quad & \sum_n \alpha_n - \frac{1}{2} \sum_{m,n} y_m y_n \alpha_m \alpha_n \phi(\mathbf{x}_m)^\top \phi(\mathbf{x}_n) \\ \text{s.t.} \quad & \sum_n \alpha_n y_n = 0 \quad \text{and} \quad 0 \leq \alpha_n \leq C, \quad \forall n \end{aligned}$$

# Summary

## Typical steps of applying Lagrangian duality

- start with a primal problem

# Summary

## Typical steps of applying Lagrangian duality

- start with a primal problem
- write down the Lagrangian (one dual variable per constraint)

# Summary

## Typical steps of applying Lagrangian duality

- start with a primal problem
- write down the Lagrangian (one dual variable per constraint)
- apply KKT conditions to find the **connections between primal and dual solutions**

# Summary

## Typical steps of applying Lagrangian duality

- start with a primal problem
- write down the Lagrangian (one dual variable per constraint)
- apply KKT conditions to find the **connections between primal and dual solutions**
- **eliminate primal variables** and arrive at the dual formulation

# Summary

## Typical steps of applying Lagrangian duality

- start with a primal problem
- write down the Lagrangian (one dual variable per constraint)
- apply KKT conditions to find the **connections between primal and dual solutions**
- **eliminate primal variables** and arrive at the dual formulation
- maximize the Lagrangian with respect to dual variables

# Summary

## Typical steps of applying Lagrangian duality

- start with a primal problem
- write down the Lagrangian (one dual variable per constraint)
- apply KKT conditions to find the **connections between primal and dual solutions**
- **eliminate primal variables** and arrive at the dual formulation
- maximize the Lagrangian with respect to dual variables
- recover the primal solutions from the dual solutions